

Introduction to Artificial Intelligence, Fall & Winter 2022
College of Computer Science, Zhejiang University
Reference Solutions for Problem Set 3

丁尧相

2023 年 1 月 5 日

Problem 2. (***) 在中心为原点的 d 维单位球 (半径为 1) 内用均匀分布采样 N 个数据点, 设至少有一个数据点落入以原点为中心, 半径为 r 的单位球的概率不低于 0.9, 此时半径 r 至少为多大?

参考解答: 可以通过计算 N 个点都没有落入球内的概率, 并令其不高于 0.1:

$$(1 - r^d)^N \leq 0.1.$$

求解可以得到

$$r \geq \sqrt[d]{1 - \sqrt[N]{0.1}}.$$

备注: 从结论中可以看出, 当维数 d 升高时, d 维球上的点越来越倾向于分布在球的外壳附近, 从而验证了维数灾难 (curse of dimensionality)。

Problem 3. 请回答下面的问题:

1. 请给出 Lec9 幻灯片第 46 页 bias&variance trade-off 的推导过程。
2. 请推导出 Lec9 幻灯片第 49 页 ridge regression 的解析解。

参考解答:

1. 参考文献 [1] 中 2.5 节推导。
2. 参考文献 [2] 中 3.4.1 节推导。

Problem 4. 在 Lec10 幻灯片的第 32 页，我们介绍了 Soft-Margin SVM 的 primal problem (原问题)，试参考前面介绍的线性可分 SVM，推导出其对应的 dual problem (对偶问题)。

参考解答:

1. 参考文献 [1] 中 6.4 节推导。

Problem 5. (***) 在介绍 AdaBoost 时，我们引入了对 SVM 中 margin 定义的推广。即对于一个二分类器 $f(\mathbf{x}) : \mathcal{X} \rightarrow [-1, +1]$ ，其在数据点 (\mathbf{x}, y) 上的 margin 定义为 $yf(\mathbf{x})$ 。试回答下面的问题：

1. 若定义多分类器 $f(\mathbf{x}) : \mathcal{X} \rightarrow \{1, 2, \dots, K\}$ ，即分类器用来预测样本 \mathbf{x} 属于 K 个类中的哪一个，你能否给出对应的 margin 的定义？
2. 理论上 AdaBoost 可以使用任何分类器作为其基学习器。在你看来，普通线性 SVM，以及使用利用核方法进行变换后的 SVM，它们作为 AdaBoost 的基学习器是否合适？为什么？

参考解答:

1. 当分类器对所有类别都有一个输出值时，可以定义为真实类别上的输出值与除真实标记外最大输出值的差。
2. 不合适。AdaBoost 是将弱分类器提升为强分类器的方法。如果弱分类器过于复杂，计算代价太高，或个体性能太强，则没有必要，甚至得不偿失。

参考文献:

- [1] 《机器学习》，周志华著。
- [2] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, “Elements of Statistical Learning”, 2nd Edition, <https://hastie.su.domains/Papers/ESLII.pdf>.